

**UNITED STATES PATENT APPLICATION FOR
SYSTEMS AND METHODS FOR STREAMING XPATH QUERY**

Inventor:

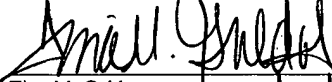
Patrick Calahan

**CERTIFICATE OF MAILING BY "EXPRESS MAIL"
UNDER 37 C.F.R. ' 1.10**

"Express Mail" mailing label number: EV 327619334

Date of Mailing: February 17, 2004

I hereby certify that this correspondence is being deposited with the United States Postal Service, utilizing the "Express Mail Post Office to Addressee" service addressed to **Mail Stop PATENT APPLICATION, Commissioner for Patents, Alexandria, VA 22313-1450**, and mailed on the above Date of Mailing with the above "Express Mail" mailing label number.


Tina M. Galdos
Signature Date: 2/17/04

SYSTEMS AND METHODS FOR STREAMING XPATH QUERY

Inventor: Patrick Calahan

COPYRIGHT NOTICE

[0001] A portion of the disclosure of this patent document contains material which is subject to copyright protection. The copyright owner has no objection to the facsimile reproduction by anyone of the patent document of the patent disclosure, as it appears in the Patent and Trademark Office patent file or records, but otherwise reserves all copyright rights whatsoever.

CLAIM OF PRIORITY

[0002] This application claims priority from the following application, which is hereby incorporated by reference in its entirety:

[0003] U.S. Provisional Application No. 60/451,313, entitled SYSTEMS AND METHODS FOR STREAMING XPATH QUERY, by Patrick Calahan, filed on February 28, 2003 (Attorney Docket No. BEAS-01330US0 SRM/DTX).

CROSS-REFERENCED CASES

[0004] The following applications are cross-referenced and incorporated herein by reference in its entirety:

[0005] U.S. Patent Application No. 10/304,207 entitled "Streaming Parser API," by Chris Fry et al., filed November 26, 2002.

FIELD OF THE INVENTION

[0006] The present invention relates to the querying of data, such as from a document or file.

BACKGROUND

[0007] XPath is a W3C language standard that can be used to address or query parts of an XML document. It models an XML document as a tree of nodes, which can include element

nodes, attribute nodes and/or text nodes. XPath can be used to identify a subset of an XML document by matching, or determining whether a node matches a pattern, similar to how SQL can be used against a database. In the typical case, an expression written in the XPath language is evaluated against an XML document to determine which parts of the document 'match' the XPath. In order to do this, the XML document must be parsed and represented in memory. One of the standard representations of XML is the Document Object Model (DOM). DOM model presents an XML document as a hierarchy of nodes through which one can navigate arbitrarily. This approach provides a lot of flexibility, but comes at a cost in terms of efficiency and memory use, as the entire document must be brought into memory at one time.

BRIEF DESCRIPTION OF THE DRAWINGS

[0008] **Figure 1** is a diagram showing an exemplary system that can be used in accordance with one embodiment of the present invention.

[0009] **Figure 2** shows an exemplary data tree that can be used with the system of Figure 1 in an embodiment.

[0010] **Figure 3** is a flowchart for an exemplary process that can be used with the system of Figure 1 in an embodiment.

DETAILED DESCRIPTION

[0011] The invention is illustrated by way of example and not by way of limitation in the figures of the accompanying drawings in which like references indicate similar elements. It should be noted that references to "an" or "one" embodiment in this disclosure are not necessarily to the same embodiment, and such references mean at least one.

[0012] Systems and methods in accordance with one embodiment of the present invention overcome deficiencies in existing XML query systems by representing the XML document as a stream of discrete 'events', with each event representing a portion of the document as the document is being parsed. Event matching can be performed against the event stream. Matching events can then be routed for processing by appropriate objects or components and returned to the event stream if necessary.

[0013] XPath can be used to identify a subset of an XML document, similar to how SQL can be used against a database. XPath is a W3C language standard that can be used to address or query parts of an XML document. It can address parts of an XML document by providing basic facilities for manipulating strings, numbers, and Boolean variables. XPath operates on the

hierarchical structure, which can be but is not limited to a tree, instead of the syntax of an XML document and can be used for matching, or determining whether a node matches a pattern. It models an XML document as a tree of nodes, which can include element nodes, attribute nodes and/or text nodes and defines a way to compute a string-value for each node type. The primary syntactic construct in XPath is the expression. An expression is evaluated to yield an object of type node-set, Boolean, number, or string. In the typical case, an expression written in the XPath language is evaluated against an XML document to determine which parts of the document 'match' the XPath. In order to do this, the XML document must be parsed and represented in memory.

[0014] Systems and methods in accordance with one embodiment of the present invention adopt a true streaming approach, passing bits of an XML document one after another, and it is up to the system to decide what to do with each bit as it passes on the stream. An advantage of a true streaming approach is that such a system is faster and far more memory efficient than a DOM-style approach, since only one portion of the document is in memory at any given time. When using a streaming parser, a system can take a stream on an XML document, generating a stream of events, one event for each node in the XML tree, and perform XPath matching on that stream. A streaming XPath system can also be schema aware, such that the system knows the XML schema for a document, that schema can be used to provide insight on how to most effectively process the document. For instance, the need to go "backwards" in a stream can be avoided if the system knows in advance which events it needs to grab and in what order those events will be received.

[0015] A streaming approach can place a greater burden on a system to maintain relevant state than a DOM approach, as a streaming approach may provide no navigation mechanisms. While such an approach provides a very efficient way to process an XML document, the efficiency comes at a cost, as there can be considerably less context available when working with a stream than when working with a DOM tree. Further, XPath has to be able to traverse the hierarchy, in some sense, in order to locate the appropriate portion of the document. In many instances, it is simple to locate an appropriate portion of XML against a DOM tree, since the system is able to walk against the tree. When using a stream, a system has to maintain context in a way that is efficient enough to make using the stream worthwhile. Some tradeoffs can be made, such as not supporting the entire XPath specification. At some point, it may be more efficient to realize an entire DOM tree, if doing a convoluted matching against the entire document.

[0016] The XPath specification defines the notion of a context, where a context is the information about an event, consisting of a node it represents, a position of the node relative to a parent node, and a function library, as well as any of several other components such as variable bindings. A location path is a type of expression that can select a set of nodes relative to the context node. The evaluation of a location path expression can result in the node-set containing the nodes being selected by the location path. Location paths can recursively contain expressions used to filter node sets. Expressions can be parsed by first dividing the character string to be parsed into tokens, then parsing the resulting token sequence.

[0017] In one embodiment, it is relatively easy to map context to the stream, as the system can maintain a stack of stream events that provide the direct ancestral line back to the root. For instance, matching an XPath that consists solely of child axes can be straightforward. In another embodiment, mapping can become more complicated in the case of descendant axes, similar to matching an entire sub-tree. In those cases, it can be necessary to spawn a tree of contexts and perform matching against each of those contexts. It can become complicated, as the system gets to maintain, and know when you can discard those cloned contexts. It can be even more complicated when matching axes called “following,” which match everything below a certain point in the document. In some cases, it is necessary to maintain that context tree and track what to add on to the tree as the system navigates its way back out of the document.

[0018] Systems and methods in accordance with one embodiment of the present invention know how to manage the multi-context mode discussed in the proceeding paragraph. They utilize the information of contexts in the stack matching against the expression to recognize when to go into this multi-context mode, when to destroy those contexts, and how to update the context stack appropriately. Certain optimizations can also be used that can know when not to match certain contexts in the context tree. XPath defines different ways to slice up a document, such as parents and children, that each has to be dealt with in a different way.

[0019] Systems and methods in accordance with one embodiment do not account for reverse axes. A reverse axis is any axis that would require going “back” through the stream. A diagram showing an exemplary “forward” and “backward” or “reverse” path through a data tree is given by **Figure 2**. A diagram of an exemplary system is shown in **Figure 1**. A streaming parser **102** generates events by parsing an XML document **100**, and then places those events on an XML event stream. Such a streaming process is demonstrated by the diagram of **Figure 3**. The streaming parser first takes a tree of an XML document as the input **300**, traverses the XML tree either through a broad-first search or a depth-first search and adds each node visited into a

data structure, e.g., a queue 302. The streaming parser then processes the queue in the first-in-first-out (FIFO) manner 304 to generate an event for the context of each node in the queue 306 and appends each event to the output stream 308. Using the event stream, the end user of the streaming API pulls events from the stream as they come through it. When a user calls for the next event on the stream, that user has a guarantee that they will get the next event. The user will find out if the next event is going to match, and will find out before the call to next returns.

[0020] In one embodiment, an XPath matching component 104 performs matching on each event received on the stream. Matching can be communicated to a caller or end user in a number of ways. These systems are doing event-based processing, as opposed to static tree-based processing. In a tree-based implementation, for example, a user can request all the nodes that match an XPath for a document. The user will receive a collection of nodes that match that XPath. Such an approach is not necessarily effective in the case of streaming, as it is then necessary to read through the document, save all the nodes, and present the collection to the user. This is fundamentally not a stream-centric way of looking at the problem. Instead, using an XPath matching approach, an observer 106 can be registered. The registered observer is an object to be notified whenever an event comes through the stream that matches this XPath. If an event matches an XPath, that event can be temporarily diverted and sent over to a user-defined object 108 that reacts to the match. Then, the event can be returned to the stream if necessary so that any subsequent object pulling events from the stream can process that event.

[0021] One embodiment may be implemented using a conventional general purpose or a specialized digital computer or microprocessor(s) programmed according to the teachings of the present disclosure, as will be apparent to those skilled in the computer art. Appropriate software coding can readily be prepared by skilled programmers based on the teachings of the present disclosure, as will be apparent to those skilled in the software art. The invention may also be implemented by the preparation of integrated circuits or by interconnecting an appropriate network of conventional component circuits, as will be readily apparent to those skilled in the art.

[0022] One embodiment includes a computer program product which is a storage medium (media) having instructions stored thereon/in which can be used to program a computer to perform any of the features presented herein. The storage medium can include, but is not limited to, any type of disk including floppy disks, optical discs, DVD, CD-ROMs, micro drive, and magneto-optical disks, ROMs, RAMs, EPROMs, EEPROMs, DRAMs, VRAMs, flash memory devices, magnetic or optical cards, nanosystems (including molecular memory ICs), or any type of media or device suitable for storing instructions and/or data.

[0023] Stored on any one of the computer readable medium (media), the present invention includes software for controlling both the hardware of the general purpose/specialized computer or microprocessor, and for enabling the computer or microprocessor to interact with a human user or other mechanism utilizing the results of the present invention. Such software may include, but is not limited to, device drivers, operating systems, execution environments/containers, and applications.

[0024] The foregoing description of the preferred embodiments of the present invention has been provided for the purposes of illustration and description. It is not intended to be exhaustive or to limit the invention to the precise forms disclosed. Many modifications and variations will be apparent to the practitioner skilled in the art. Embodiments were chosen and described in order to best describe the principles of the invention and its practical application, thereby enabling others skilled in the art to understand the invention, the various embodiments and with various modifications that are suited to the particular use contemplated. It is intended that the scope of the invention be defined by the following claims and their equivalents.